1·0

2·8    2·5

3·15   2·2

3·5

1·1    4·0    2·0

4·5    1·8

1·25   1·4   1·6

NATIONAL BUREAU OF STANDARDS
MICROCOPY RESOLUTION TEST CHART

# COMPUTER SCIENCE
# TECHNICAL REPORT SERIES

# UNIVERSITY OF MARYLAND
## COLLEGE PARK, MARYLAND
### 20742

TR-571

12 10 p.

# Sensitivity Coefficients for the Effects of Errors in the Independent Variables in a Linear Regression.

by

G. W. Stewart

D D C

OCT 4 1977

C

## Abstract

This paper is concerned with errors in the observed values of the independent variables of a linear regression. We propose sensitivity coefficients to measure the effects of these errors and show that they can easily be computed from quantities ordinarily calculated in performing the regression.

409 022

# Sensitivity Coefficients
## for the Effects of Errors in
## the Independent Variables in
## a Linear Regression
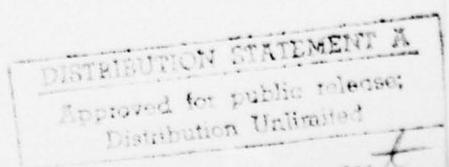
G. W. Stewart

## 1. Introduction

In this paper we shall be concerned with the regression problem

$$\text{minimize } \|y - X\beta\|^2 ,$$

where $X$ is an $n \times p$ matrix of rank $p$, $y$ is an n-vector, and $\|\cdot\|$ denotes the usual Euclidean vector norm.* The problem has the unique solution

$$(1.1) \qquad \beta = X^{\dagger} y$$

where $X^{\dagger} = (X'X)^{-1} X^{T}$ is the pseudo-inverse of $X$.

Although classical regression theory concerns itself with the statistical analysis of errors in the vector y, it frequently happens that the design matrix $X$ is itself contaminated with errors, so that one is effectively working with a perturbed matrix $X + E$. For example, the columns of $X$ may be measured by means of some instrument for which the originator of the problem can only give crude error estimates. In this case the data analyst is faced with the problem of deciding when the effects of the errors can be ignored. The problem is especially

---

* In the sequel $\|\cdot\|$ will also denote the spectral matrix norm defined by $\|A\| = \sup \{\|Ax\| : \|x\| = 1\}$. See [3] for details.

critical in general purpose regression routines, where it is desirable
to provide the user with a set of easily interpretable numbers that
indicate the magnitude of the effects of the errors.

A partial solution is provided by the perturbation theory for the
regression problem (for a survey of this theory see [4]). This theory
bounds perturbations in $\beta$ in terms of $\|E\|$ and of the "condition
number" $\kappa = \|X\|\|X^{\dagger}\|$. Although the results of this theory shed considerable
light on the behavior of regression problems under perturbations in $X$,
they are unsatisfactory in practice for two reasons. First they bound
only the norm of the perturbation in $\beta$, so that a large perturbation in
one component can conceal the fact that the others have small perturba-
tions. More important, they are not scale invariant; changing the scale
of the columns of $X$ will change $\kappa$, even though the statistical problem
is essentially unaltered. This phenomenon makes the results quite diffi-
cult to interpret.

Taking a different approach, Beaton, Rubin, and Barone [1] have
derived measures of sensitivity that to some extent answer the above ob-
jections. However, it is assumed that the errors are unbiased and $n$ is
large. Swindel and Bauer [5] have derived a useful bound for the relative
bias in $\beta$, which measures the relative effects of perturbations in $X$
compared to the usual statistical errors in $y$. In this paper we shall
derive coefficients $\gamma_{ij}$ that measure the sensitivity of $\beta_i$ to changes
in column $j$ of $X$. Specifically $\gamma_{ij}$ is the norm of the Frechet

derivative of $\beta_i$ regarded as a function of the j-th column of $\underset{\sim}{X}$.
If $\varepsilon$ is the norm of the perturbation of the i-th column of $\underset{\sim}{X}$, then
$\gamma_{ij}\varepsilon$ will be an asymptotic bound on the perturbation induced in $\beta_j$.*

## 2. Derivation of the Coefficients

Although it is in principle possible to calculate the required
derivatives directly from the normal equations $(\underset{\sim}{X}'\underset{\sim}{X})\beta = \underset{\sim}{X}'y$, we prefer
to approach the problem through a first order perturbation theorem that
is useful in its own right.

Theorem 2.1. In the notation of the last section, let $\underset{\sim}{\beta}$ be the
solution of the regression problem (1.1), and let $\underset{\sim}{r} = \underset{\sim}{y} - \underset{\sim}{X}\underset{\sim}{\beta}$ be the
corresponding residual vector. Let $\underset{\sim}{E}$ be an $n \times p$ matrix. If

$$(2.1) \qquad \qquad \|\underset{\sim}{X}^{+}\underset{\sim}{E}\| < 1 ,$$

then $\underset{\sim}{X} + \underset{\sim}{E}$ has rank $p$ so that there is a unique solution $\underset{\sim}{\tilde{\beta}}$ of the
problem

$$\text{minimize } \|\underset{\sim}{y} - (\underset{\sim}{X}+\underset{\sim}{E})\underset{\sim}{\tilde{\beta}}\|^2 .$$

Moreover, as $\underset{\sim}{E}$ approaches zero

$$\underset{\sim}{\tilde{\beta}} = \underset{\sim}{\beta} - \underset{\sim}{X}\,\underset{\sim}{E}\beta + (\underset{\sim}{X}'\underset{\sim}{X})^{-1}\underset{\sim}{E}^{T}\underset{\sim}{r} + O(\|\underset{\sim}{E}\|^2) .$$

Proof. For a proof that (2.1) implies that $\underset{\sim}{X} + \underset{\sim}{E}$ has full rank, see
[4]. This implies that for all sufficiently small $\underset{\sim}{E}$, $\underset{\sim}{\tilde{\beta}}$ exists and is
given by

---

* Here, and throughout this note, the term asymptotic refers to behavior for
small $\varepsilon$, not large $n$.

$$\tilde{\beta} = [(X+E)'(X+E)]^{-1}(X+E)y \ .$$

Now it is well known [3] that if $P$ is sufficiently small, then $I + P$ is nonsingular and

$$(I+P)^{-1} = I - P + O(\|P\|^2) \ .$$

It follows that

$$[(X+E)'(X+E)]^{-1} = (X'X+X'E+E'X+E'E)^{-1}$$

$$= \{X'X[I+(X'X)^{-1}(X'E+E'X+E'E)]\}^{-1}$$

$$= [I-(X'X)^{-1}(X'E+E'X)](X'X)^{-1} + O(\|E\|^2) \ .$$

Hence

$$\tilde{\beta} = (X'X)^{-1}X'y + (X'X)^{-1}E'y - (X'X)^{-1}X'E(X'X)^{-1}X'y$$

$$- (X'X)^{-1}E'X(X'X)^{-1}X'y + O(\|E\|^2)$$

$$= \beta - X^{\dagger}E\beta + (X'X)^{-1}E'(y-X\beta)$$

$$= \beta - X^{\dagger}E\beta + (X'X)^{-1}E^{T}r \ . \ \square$$

Theorem 2.1 immediately gives an expression for the Frechet derivative of $\beta_i$ regarded as a function of the j-th column of X.

Corollary 2.2. Let $\beta_i = f_{ij}(x_j)$, where $x_j$ denotes the j-th column of X. Then

$$(2.2) \qquad Df_{ij} = -\beta_j e_i^T X^{\dagger} + e_i'(X'X)^{-1}e_j r^T \ .$$

$\underline{\text{Proof}}$. The Frechet derivative $\underset{\sim}{D}f_{ij}$ is the unique row vector satisfying

(2.3) $\qquad f_{ij}(\underset{\sim}{x}_j + \underset{\sim}{e}) = \beta_i + \underset{\sim}{D}f_{ij}\underset{\sim}{e} + O(\|\underset{\sim}{e}\|^2)$ .

Let $\underset{\sim}{e}_j$ denote the j-th unit vector. Then a perturbation $\underset{\sim}{e}$ in $\underset{\sim}{x}_j$ amounts to adding to $\underset{\sim}{X}$ the matrix $\underset{\sim}{E} = \underset{\sim}{e}\underset{\sim}{e}_j^T$. Hence from Theorem 2.1

$$f_{ij}(\underset{\sim}{x}_j + \underset{\sim}{e}) = \underset{\sim}{e}_i'[\underset{\sim}{\beta} - \underset{\sim}{X}(\underset{\sim}{e}\underset{\sim}{e}_j')\underset{\sim}{\beta} + (\underset{\sim}{X}'\underset{\sim}{X})^{-1}(\underset{\sim}{e}\underset{\sim}{e}_j')'\underset{\sim}{r}] + O(\|\underset{\sim}{e}\|^2)$$

$$= \beta_i - \beta_j \underset{\sim}{e}_i'\underset{\sim}{X}\, \underset{\sim}{e} + \underset{\sim}{e}_i'(\underset{\sim}{X}\underset{\sim}{X})^{-1}\underset{\sim}{e}_j\underset{\sim}{e}'\underset{\sim}{r} + O(\|\underset{\sim}{e}\|^2)$$

$$= \underset{\sim}{\beta}_i - (\beta_j\underset{\sim}{e}_i'\underset{\sim}{X} - \underset{\sim}{e}_i'(\underset{\sim}{X}'\underset{\sim}{X})^{-1}\underset{\sim}{e}_j\underset{\sim}{r}')\underset{\sim}{e} + O(\|\underset{\sim}{e}\|^2) \, ,$$

which shows that $f_{ij}'$ defined by (2.2) satisfies (2.3).□

$\underline{\text{Corollary}}$ 2.3. Let $C = (\underset{\sim}{X}'\underset{\sim}{X})^{-1}$. Then

$$\|\underset{\sim}{D}f_{ij}\| = \gamma_{ij} \equiv \sqrt{|\beta_j|^2 c_{ii} + \|r\|^2 c_{ij}^2} \, .$$

$\underline{\text{Proof}}$. The vector $r$ is orthogonal to the column space of $X$, which is the same as the row space of $X^\dagger$. Hence $e_i^T X^\dagger$ and $r'$ are orthogonal, and

$$\|Df_{ij}\|^2 = |\beta_j|^2 \|e_i X^\dagger\|^2 + c_{ij}^2 \|r\|^2 \, .$$

The proof will be complete if we can show that $\|e_i' X^\dagger\|^2 = c_{ii}$. But

$$\|e_i' X^\dagger\|^2 = \|e_i'(X'X)^{-1}X'\| = e_i'(X'X)^{-1}X'X(X'X)^{-1}e_i = e_i'(X'X)^{-1}e_i = c_{ii} \, . \square$$

## 3. Applications

It should be observed that the numbers $\gamma_{ij}$ can be calculated from the quantities that are usually computed in the course of the regression. For example, if sweep techniques (e.g. see [2]) have been used to solve the problem, the numbers $c_{ij}$ will be available, as will $\|\underset{\sim}{r}\|^2$, since it is simply the residual sum of squares.

The results are readily interpretable. If a bound $\varepsilon_j$ on the norm of the error in the j-th column is available, then the error $\delta_i$ induced in $\beta_i$ is asymptotically bounded by $\gamma_{ij}\varepsilon_j$. If this number is too large, the problem requires further study. In interpreting the results it should always be borne in mind that $\gamma_{ij}\varepsilon_j$ is a first order bound. A large value is a signal that something may be wrong, but if $\varepsilon_j$ is so large that the first order approximation is not applicable, then the difficulties may turn out to be illusory.

A particularly attractive feature of the asymptotic bounds is that, since they deal with indivual components of $\beta$, their interpretation is independent of the scaling of the columns of $\underset{\sim}{X}$. This is particularly apparent if the bounds are cast in terms of relative error in the form

$$\frac{|\delta_i|}{|\beta_i|} \lesssim \left( \frac{\gamma_{ij}\|\underset{\sim}{x}_j\|}{\|\beta_i\|} \right) \frac{\varepsilon_j}{\|\underset{\sim}{x}_j\|} \quad .$$

Now

$$(3.1) \qquad \frac{\gamma_{ij}\|\underset{\sim}{x}_j\|}{|\beta_i|} = \left[ (|\beta_j|^2\|\underset{\sim}{x}_j\|^2)\left(\frac{c_{ii}}{|\beta_i|^2}\right) + (\|r\|^2)\left(\frac{|c_{ij}|^2\|\underset{\sim}{x}_j\|^2}{|\beta_i|^2 \ \|}\right) \right]^{\frac{1}{2}} ,$$

and each parenthesized term in the right hand side of (3.1) is easily seen to be invariant under scaling of the columns of $\underset{\sim}{X}$. Indeed, in some applications where the $\beta_i$'s are known to be bounded away from zero it may be more appropriate to report $\gamma_{ij} \|\underset{\sim}{x}_j\| / |\beta_i|$ than $\gamma_{ij}$.

In deriving our bounds, we have used the Cauchy-Schwarz inequality $\|\underset{\sim}{x}^*\underset{\sim}{y}\| \leq \|\underset{\sim}{x}\| \|\underset{\sim}{y}\|$, an inequality which is usually pessimistic, since it must account for the worst case where $\underset{\sim}{x}$ and $\underset{\sim}{y}$ are dependent. If we are willing to assume more about the perturbation $\underset{\sim}{e}$ in $\underset{\sim}{x}_j$, then we may be able to say more. For example, we have the following consequence of Corollary 2.3.

**Corollary** 3.1. Let $\underset{\sim}{e} \in N(0, \sigma^2 \underset{\sim}{I})$. Then $\underset{\sim}{D} f_{ij} \underset{\sim}{e}$ is normally distributed with mean zero and standard deviation $\gamma_{ij} \sigma$.

## References

1. A. E. Beaton, D. B. Rubin, and J. L. Barne, The acceptability of regression solutions: another look at computational accuracy, J. Amer. Stat. Assoc. 71 (1976) 158-168.

2. A. P. Dempster, Elements of Continuous Multivariate Analysis, Addison-Wesley, Reading, Massachusetts (1969).

3. G. W. Stewart, Introduction to Matrix Computations, Academic Press, New York (1973).

4. _____, On the perturbation of pseudo-inverses, projections, and linear least squares problems, to appear SIAM Rev.

5. B. F. Swindel and D. R. Bower, Rounding errors in the independent variables in a general linear model, Technometrics 14 (1972) 215-218.

SECURITY CLASSIFICATION OF THIS PAGE *(When Data Entered)*

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| **1. REPORT NUMBER** <br> ONR-N00014-76-C-0391-571 | **2. GOVT ACCESSION NO.** | **3. RECIPIENT'S CATALOG NUMBER** |
| **4. TITLE** *(and Subtitle)* <br> SENSITIVITY COEFFICIENTS FOR THE EFFECTS OF ERRORS IN THE INDEPENDENT VARIABLES IN A LINEAR REGRESSION | | **5. TYPE OF REPORT & PERIOD COVERED** <br> Technical Report |
| | | **6. PERFORMING ORG. REPORT NUMBER** <br> TR-571 |
| **7. AUTHOR(s)** <br> G. W. Stewart | | **8. CONTRACT OR GRANT NUMBER(s)** <br> N00014-76-C-0391 |
| **9. PERFORMING ORGANIZATION NAME AND ADDRESS** <br> Department of Computer Science <br> University of Maryland <br> Colege Park, MD 20742 | | **10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS** |
| **11. CONTROLLING OFFICE NAME AND ADDRESS** <br> Mathematics Branch <br> Office of Naval Research <br> Arlington, VA 22217 | | **12. REPORT DATE** <br> September 1977 |
| | | **13. NUMBER OF PAGES** <br> 9 |
| **14. MONITORING AGENCY NAME & ADDRESS***(if different from Controlling Office)* | | **15. SECURITY CLASS.** *(of this report)* <br> UNCLASSIFIED |
| | | **15a. DECLASSIFICATION/DOWNGRADING SCHEDULE** |

**16. DISTRIBUTION STATEMENT** *(of this Report)*

Approved for public release; distribution unlimited.

**17. DISTRIBUTION STATEMENT** *(of the abstract entered in Block 20, if different from Report)*

**18. SUPPLEMENTARY NOTES**

**19. KEY WORDS** *(Continue on reverse side if necessary and identify by block number)*

regression
least squares
errors in variables
sensitivity

**20. ABSTRACT** *(Continue on reverse side if necessary and identify by block number)*

This paper is concerned with errors in the observed values of the independent variables of a linear regression. We propose sensitivity coefficients to measure the effects of these errors and show that they can easily be computed from quantities ordinarily calculated in performing the regression.

**DD** <sub></sub> FORM 1 JAN 73 **1473** EDITION OF 1 NOV 65 IS OBSOLETE

SECURITY CLASSIFICATION OF THIS PAGE *(When Data Entered)*